

wCorr Arguments

Paul Bailey

2016-03-03

This vignette explores two Boolean switches in the package. First, the `ML` switch allows for either a non-MLE (but consistent) estimate of the nuisance parameters that define the binning process to be used (`ML=FALSE`) or for the nuisance parameters to be estimated using the MLE (`ML=TRUE`). Second the `fast` argument gives the option to use a pure R implementation (`fast=FALSE`) or an implementation that relies on the `Rcpp` and `RcppArmadillo` packages (`fast=TRUE`).

Numerical simulations show that the results are essentially unaffected by either of these switches and so it is recommended to use `fast=TRUE` and `ML=FALSE` which will drastically speed computation.

The “wCorr Formulas” vignette describes the statistical properties of the correlation estimators in the package.

ML switch

The correlation coefficients between two vectors of random variables that are jointly bivariate normal—call the vectors \mathbf{X} and \mathbf{Y} .

$$\begin{pmatrix} X \\ Y \end{pmatrix} \sim N \left[\begin{pmatrix} \mu_x \\ \mu_y \end{pmatrix}, \Sigma \right]$$

where $N(\mathbf{A}, \Sigma)$ is the bivariate normal distribution with mean \mathbf{A} and covariance Σ .

The i^{th} members of the vectors are then called x_i and y_i .

Polyserial computation

The derivatives of ℓ can be computed but are not readily computed and so when the `ML` argument is set to `FALSE` (the default) a one dimensional optimization of ρ is calculated using `stats::optimize`. When the `ML` argument is set to `TRUE` a multi-dimensional optimization is done for ρ and $\boldsymbol{\theta}$ using `minqa::bobyqa`. As is shown below, the difference between these two is slight, if present, and so the default value of `ML` is recommended.

Because the optimization is not perfect when the correlation is in a boundary condition ($\rho \in \{-1, 1\}$), a check for perfect correlation is performed before the above optimization by simply seeing if the values of \mathbf{X} and \mathbf{M} have exactly the same order.

Polychoric computation

This again mirrors the treatment of the polyserial. The derivatives of ℓ can be computed but are not readily computed and so when the `ML` argument is set to `FALSE` (the default) a one dimensional optimization of ρ is calculated using `stats::optimize`. When the `ML` argument is set to `TRUE` a multi-dimensional optimization is done for ρ , $\boldsymbol{\theta}$, and $\boldsymbol{\theta}'$ using `minqa::bobyqa`. As is shown below, the difference between these two is slight, if present, and so the default value of `ML` is recommended.

Because the optimization is not perfect when the correlation is in a boundary condition ($\rho \in \{-1, 1\}$), a check for perfect correlation is performed before the above optimization by simply seeing if the values of \mathbf{P} and \mathbf{M} have a Goodman-Kruskal correlation coefficient of -1 or 1. When this is the case, the MLE of -1 or 1, respectively, is returned.

General setup for the unweighted case

A simulation is run several times using the following method:

- selecting a value of n (the number of observations)
- selecting a true correlation coefficient ρ
- generating \mathbf{X} and \mathbf{Y}
- selecting the value of t and t' (the number of bins for \mathbf{M} and \mathbf{P})
- selecting θ and θ'
- confirming that at least 2 levels of \mathbf{M} and \mathbf{P} are occupied (if not, re-run to generating \mathbf{X} and \mathbf{Y})
- calculating and recording relevant statistics

ML switch

It is easy to prove the consistency of the θ for the polyserial and θ and θ' using the non-ML case. Similarly, for ρ , because it is an MLE that can be obtained by taking a derivative and setting it equal to zero, the results are asymptotically unbiased and obtain the Cramer-Rao lower bound.

This does not speak to the small sample properties of these correlation coefficients. Previous work has described their properties by simulation and so that tradition is continued below.

- plot that shows difference as a function of ρ and at $n = 10$ and $n = 1000000$ between `ML=FALSE` and `ML=TRUE` when `fast=TRUE`

fast switch

This section looks at the agreement between the pure R implementation of the optimizations and the `Rcpp` and `RcppArmadillo` implementation. The code can compute with either option by setting `fast=FALSE` (pure R) or `fast=TRUE` (`Rcpp`).

This is the summary of all differences between the `fast=TRUE` and `fast=FALSE` runs for the polyserial

- plot that shows difference as a function of ρ and at $n = 10$ and $n = 1000000$ between `fast=FALSE` and `fast=TRUE` when `ML=TRUE`
- plot that shows difference as a function of ρ and at $n = 10$ and $n = 10000$ between `fast=FALSE` and `fast=TRUE` when `ML=FALSE`

Implications for speed

- plot of compute time vs n for all four options. Each stop when the mean ≥ 20 seconds